

Redundancy of Exchangeable Estimators

Narayana P. Santhanam¹, Anand D. Sarwate² and Jae Oh Woo³

October 22, 2014

Abstract

Exchangeable random partition processes are the basis for Bayesian approaches to statistical inference in large alphabet settings. On the other hand, the notion of the *pattern* of a sequence provides an information-theoretic framework for data compression in large alphabet scenarios. Because data compression and parameter estimation are intimately related, we study the redundancy of Bayes estimators coming from Poisson-Dirichlet priors (or “Chinese restaurant processes”) and the Pitman-Yor prior. This provides an understanding of these estimators in the setting of unknown discrete alphabets from the perspective of universal compression. In particular, we identify relations between alphabet sizes and sample sizes where the redundancy is small, thereby characterizing useful regimes for these estimators.

Keywords— exchangeability, random exchangeable partitions, Chinese restaurant process, Pitman Yor process, strong and weak universal compression

1 Introduction

A number of statistical inference problems of significant contemporary interest, such as text classification, language modeling and DNA microarray analysis, require inferences based on observed sequences of symbols in which the sequence length or sample size is comparable or even smaller than the set of symbols, the alphabet. For instance, language models for speech recognition estimate distributions over English words using text examples much smaller than the vocabulary.

Inference in this setting has received a lot of attention, from Laplace [1, 2, 3] in the 18th century, to Good [4] in the mid-20th century, to an explosion of work in the statistics [5, 6, 7, 8, 9, 10, 11, 12, 13], information theory [14, 15, 16, 17, 18, 19] and machine learning [20, 21, 22, 23] communities in the last few decades. A major strand in the information theory literature on the subject has been based on the notion of patterns. The pattern of a sequence characterizes the repeat structure in the sequence, which is the information that can be described well (see Orlitsky *et al.* [24] for formal characterizations of this idea). The statistical literature has emphasized the importance of exchangeability, which generalizes the notion of independence.

¹Department of Electrical Engineering, University of Hawaii at Manoa, 2540 Dole Street, Honolulu, HI 96822, USA, nsanthan@hawaii.edu

²Department of Electrical and Computer Engineering, Rutgers, The State University of New Jersey, 94 Brett Road, Piscataway, NJ 08854, USA, asarwate@ece.rutgers.edu

³Applied Mathematics Program, Yale University, 51 Prospect St, New Haven, CT 06511, USA, jaeoh.woo@yale.edu

We consider measures over infinite sequences $X_1, X_2 \dots$, where X_i come from a countable (infinite) set (the alphabet). Let \mathcal{I} be the collection of all distributions over countable (potentially infinite) alphabets. For $p \in \mathcal{I}$, let $p^{(n)}$ denote the product distribution corresponding to an independent and identically distributed (i.i.d.) sample $X_1^n = (X_1, X_2, \dots, X_n)$, where $X_i \sim p$. Let $\mathcal{I}^{(n)}$ be the collection of all such i.i.d. distributions on length n sequences drawn from countable alphabets. Let \mathcal{I}^∞ be the collection of all measures over infinite sequences of symbols from countable alphabets $X_1, X_2 \dots$, where the $\{X_i\}$ are i.i.d. according to some distribution in \mathcal{I} . The measures are constructed by extending to the Borel sigma algebra the i.i.d. probability assignments on finite length sequences, namely $\mathcal{I}^{(n)}$, $n \geq 1$. We call \mathcal{I}^∞ the set of i.i.d. measures.

Based on a sample $X_1^n = (X_1, X_2, \dots, X_n)$ from an unknown $p^{(n)} \in \mathcal{I}^{(n)}$ (or equivalently, the corresponding measure in \mathcal{I}^∞), we want to create an estimator q_n , which assigns probabilities to length- n sequences. We are interested in the behavior of the sequence of estimators $\{q_n : n = 1, 2, \dots\}$. With some abuse of notation, we will use q to denote the estimator q_n when the sample size n is clear from context. We want q_n to approximate $p^{(n)}$ well; in particular, we would like q_n to neither overestimate nor underestimate the probability of sequences of length n under the true $p^{(n)}$.

Suppose that there exist $R_n > 0$ and $A_n > 0$, such that for each $p \in \mathcal{I}$, we have:

$$p^{(n)}\left(\left\{X_1^n : q_n(X_1^n) > R_n p^{(n)}(X_1^n)\right\}\right) < 1/A_n.$$

If A_n is any function of n that grows sufficiently fast with n , the sequence $\{q_n\}$ does not asymptotically overestimate probabilities of length- n sequences by a factor larger than R_n with probability one, no matter what measure $p \in \mathcal{I}$ generated the sequences.

Protecting against underestimation is not so simple. The redundancy of an estimator q_n (defined formally in Section 2.4) for a length n sequence x_1^n measures how closely $q_n(x_1^n)$ matches:

$$\max_{p^{(n)} \in \mathcal{I}^{(n)}} p^{(n)}(x_1^n),$$

the largest probability assigned to x_1^n by any distribution in $\mathcal{I}^{(n)}$. The estimator redundancy usually either maximizes the redundancy of a sequence or takes the expectation over all sequences. Ideally, we want the estimator redundancy to grow sublinearly in the sequence length n , so that the per-sample redundancy vanishes as $n \rightarrow \infty$. If so, we call the estimator universal for \mathcal{I} . Redundancy thus captures how well q performs against the collection \mathcal{I} , but the connections between estimation problems and compression run deeper.

In this paper, we consider estimators formed by taking a measure (prior) on \mathcal{I} . Different priors induce different distributions on the data X_1^n . We think of the prior as randomly choosing a distribution p in \mathcal{I} , and the observed data X_1^n is generated according to this p . How much information about the underlying distribution p can we obtain from the data (assuming we know the prior)? Indeed, a well known result [25, 26, 27] proves that the redundancy of the best possible estimator for \mathcal{I}^∞ equals the maximum (over all choices of priors) information (in bits) that is present about the underlying source in a length n sequence generated in this manner.

Redundancy is well defined for finite alphabets; recent work [16] has formalized a similar framework for countably infinite alphabets. This framework is based on the notion of patterns of sequences that abstract the identities of symbols and indicate only the relative order of appearance. For example, the pattern of FEDERER is 1232424, while that of PATTERN is 1233456. The crux of the idea is that instead of considering the set of measures \mathcal{I}^∞ over infinite sequences, we consider the set of measures induced over patterns of the sequences. It then follows that now our estimate q_ψ is a measure over patterns. While the variables in the sequence are i.i.d., the

corresponding pattern merely corresponds to a exchangeable random partition. We can associate a predictive distribution with the pattern probability estimator q_Ψ . This is an estimate of the distribution of X_{n+1} given the previous observations, and it assigns probabilities to the event that X_{n+1} will be “new” (has not appeared in X_1^n) and probabilities to the events that X_{n+1} takes on one of the values that has been seen so far.

The above view of estimation also appears in the statistical literature on Bayesian nonparametrics that focuses on exchangeability. Kingman [28] advocated the use of exchangeable random partitions to accommodate the analysis of data from an alphabet that is not bounded or known in advance. A more detailed discussion of the history and philosophy of this problem can be found in the works of Zabell [11, 29] collected in [30]. One of the most popular exchangeable random partition processes is the “Chinese restaurant process” [10], which is a special case of the Poisson–Dirichlet or Pitman–Yor process [31, 13]. These processes can be viewed as prior distributions on the set of all discrete distributions that can be used as the basis for estimating probabilities and computing predictive distributions.

In this paper, we evaluate the performance of the sequential estimators corresponding to these exchangeable partition processes. As described before, \mathcal{I} is the collection of all distributions over countable (potentially infinite) alphabets, and \mathcal{I}^∞ is the collection of all i.i.d. measures with single letter marginals in \mathcal{I} . Let \mathcal{I}_Ψ be the collection of all measures over patterns induced by measures in \mathcal{I}^∞ . We evaluate the redundancy of estimators based on the Chinese restaurant process (CRP), the Pitman–Yor (PY) process and the Ewen’s sampling formula against \mathcal{I}_Ψ .

In the context of sequential estimation, early work [16] showed that for the collection \mathcal{I}_Ψ of measures over patterns, universal estimators do exist: the normalized redundancy is $O(n^{1/2})$. More recent work [32, 33] proved tight bounds on worst-case and average redundancy; these results show that there are sequential estimators with normalized redundancy of $O(n^{1/3})$. However, these estimators are computationally intensive and (generally speaking) infeasible in practice. Acharya et al. [34] demonstrated a linear-time estimator with average redundancy $O(n^{1/2})$, improving over the earlier constructions achieving $O(n^{2/3})$ [16]. By contrast, estimators such as the CRP or the PY estimator were not developed in a universal compression framework, but have been very successful from a practical standpoint. The goal of this paper is to understand these Bayesian estimators from the universal compression perspective.

For the case of the estimators studied in nonparametric Bayesian statistics, our results show that they are in general neither weakly nor strongly universal when compressing patterns or equivalently exchangeable random partitions. While the notion of redundancy is in some sense different from other measures of accuracy, such as the concentration of the posterior distribution about the true distribution, the parameters of the CRPs or the PY processes that do compress well often correspond to the maximum likelihood estimates obtained from the sample.

Because we choose to measure redundancy in the worst case over p , the underlying alphabet size may be arbitrarily large with respect to the sample size n . Consequently, for a fixed sample of size n , the number of symbols could be large, for example a constant fraction of n . The CRP and PY estimators do not have good redundancy against such samples, since they are not the cases the estimators are designed for. However, we can show that a mixture of estimators corresponding to CRP estimators is weakly universal. This mixture is made by optimizing individual CRP estimators that (implicitly) assume a bound on the support of p . If such a bound is known in advance, we can derive much tighter bounds on the redundancy. In this setting, the two-parameter Poisson–Dirichlet (or Pitman–Yor) estimator is superior to the estimator derived from the Chinese restaurant process.

In order to describe our results, we require a variety of definitions from different research

communities. In the next section, we describe this preliminary material and place it in context before describing the main results in Section 3.

2 Preliminaries

In this paper, we use the “big- O ” notation. A function $f(n) = O(g(n))$ if there exists a positive constant C , such that for sufficiently large n , $|f(n)| \leq C|g(n)|$. A function $f(n) = \Omega(g(n))$, if there exists a positive constant C' , such that for sufficiently large n , $|f(n)| \geq C'|g(n)|$. A function $f(n) = \Theta(g(n))$, if $f(n) = O(g(n))$ and $f(n) = \Omega(g(n))$.

Let \mathcal{I}_k denote the set of all probability distributions on alphabets of size k and \mathcal{I}_∞ be all probability distributions on countably infinite alphabets, and let:

$$\mathcal{I} = \mathcal{I}_\infty \cup \bigcup_{k \geq 1} \mathcal{I}_k$$

be the set of all discrete distributions irrespective of support and support size.

For a fixed p , let $x_1^n = (x_1, x_2, \dots, x_n)$ be a sequence drawn i.i.d. according to p . We denote the pattern of x_1^n by ψ_1^n . The pattern is formed by taking $\psi_1 = 1$ and:

$$\psi_i = \begin{cases} \psi_j & x_i = x_j, j < i \\ 1 + \max_{j < i} \psi_j & x_i \neq x_j, \forall j < i \end{cases}$$

For example, the pattern of $x_1^7 = \text{FEDERER}$ is $\psi_1^7 = 1232424$. Let ψ^n be the set of all patterns of length n . We write $p(\psi^n)$ for the probability that a length- n sequence generated by p has pattern ψ^n . For a pattern ψ_1^n , we write ϕ_μ for the number of symbols that appear μ times in ψ_1^n , and $m = \sum \phi_\mu$ is the number of distinct symbols in ψ_1^n . We call ϕ_μ the prevalence of μ . Thus, for FEDERER, we have $\phi_1 = 2$, $\phi_2 = 1$, $\phi_3 = 1$ and $m = 4$.

2.1 Exchangeable Partition Processes

An exchangeable random partition refers to a sequence $(C_n : n \in \mathbb{N})$, where C_n is a random partition of the set $[n] = \{1, 2, \dots, n\}$, satisfying the following conditions: (i) the probability that C_n is a particular partition depends only on the vector (s_1, s_2, \dots, s_n) , where s_k is the number of parts in the partition of size k ; and (ii) the realizations of the sequence are consistent in that all of the parts of C_n are also parts of the partition C_{n+1} , except that the new element $n + 1$ may either be in a new part of C_{n+1} by itself or has joined one of the existing parts of C_n .

For a sequence X_1, \dots, X_n from a discrete alphabet, one can partition the set $[n]$ into component sets $\{A_x\}$, where $A_x = \{i : X_i = x\}$ are the indices corresponding to the positions in which x has appeared. When the $\{X_i\}$ are drawn i.i.d. from a distribution in \mathcal{I} , the corresponding sequence of random partitions is called a paintbox process.

The remarkable Kingman representation theorem [8] states that the probability measure induced by any exchangeable random partition is a mixture of paintbox processes, where the mixture is taken using a probability measure (“prior” in Bayesian terminology) on the class of paintbox processes. Since each paintbox process corresponds to a discrete probability measure (the one such that i.i.d. X_i drawn from it produced the paintbox process), the prior may be viewed as a distribution on the set of probability measures on a countable alphabet. For technical reasons, the alphabet is assumed to be hybrid, with a discrete part, as well as a continuous part, and also, one needs to work with the space of ordered probability vectors (see [35] for details).

2.2 Dirichlet Priors and Chinese Restaurant Processes

Not surprisingly, special classes of priors give rise to special classes of exchangeable random partitions. One particularly nice class of priors on the set of probability measures on a countable alphabet is that of the Poisson–Dirichlet priors [36, 5, 37] (sometimes called Dirichlet processes, since they live on the infinite-dimensional space of probability measures and generalize the usual finite-dimensional Dirichlet distribution).

The Chinese restaurant process (or CRP) is related to the so-called Griffiths–Engen–McCloskey (GEM) distribution with parameter θ , denoted by $\text{GEM}(\theta)$. Consider W_1, W_2, \dots drawn i.i.d. according to a $\text{Beta}(1, \theta)$ distribution, and set:

$$\begin{aligned} p_1 &= W_1 \\ p_i &= W_i \prod_{j < i} (1 - W_j) \quad \forall i > 1 \end{aligned}$$

This can be interpreted as follows: take a stick of unit length and break it into pieces of size W_1 and $1 - W_1$. Now take the piece of size $1 - W_1$ and break off a W_2 fraction of that. Continue in this way. The resulting lengths of the sticks create a distribution on a countably infinite set. The distribution of the sequence $p = (p_1, p_2, \dots)$ is the $\text{GEM}(\theta)$ distribution.

Remark Let π denote the elements of p sorted in decreasing order, so that $\pi_1 \geq \pi_2 \geq \dots$. Then, the distribution of π is the Poisson–Dirichlet distribution $\text{PD}(\theta)$ as defined by Kingman. \square

Another popular class of distributions on probability vectors is the Pitman–Yor family of distributions [13], also known as the two-parameter Poisson–Dirichlet family of distributions $\text{PD}(\alpha, \theta)$. The two parameters here are a discount parameter $\alpha \in [0, 1]$ and a strength parameter $\theta > -\alpha$. The distribution $\text{PD}(\alpha, \theta)$ can be generated in a similar way as the Poisson–Dirichlet distribution $\text{PD}(\theta) = \text{PD}(0, \theta)$ described earlier. Let each W_i be drawn independently according to a $\text{Beta}(1 - \alpha, \theta + i\alpha)$ distribution, and again, set:

$$\begin{aligned} \tilde{p}_1 &= W_1 \\ \tilde{p}_i &= W_i \prod_{j < i} (1 - W_j) \quad \forall i > 1 \end{aligned}$$

A similar “stick-breaking” interpretation holds here, as well. Now, let p be equal to the sequence \tilde{p} sorted in descending order. The distribution of p is $\text{PD}(\alpha, \theta)$. If we have $\alpha < 0$ and $\theta = r|\alpha|$ for integer r , we may obtain a symmetric Dirichlet distribution of dimension r .

2.3 Pattern Probability Estimators

Given a sample x_1^n with pattern ψ^n , we would like to produce a pattern probability estimator. This is a function of the form $q(\psi_{n+1} | \psi^n)$ that assigns a probability of seeing a symbol previously seen in ψ^n , as well as a probability of seeing a new symbol. In this paper, we will investigate two different pattern probability estimators based on Bayesian models.

The Ewens sampling formula [38, 39, 40], which has its origins in theoretical population genetics, is a formula for the probability mass function of a marginal of a CRP corresponding to a fixed population size. In other words, it specifies the probability of an exchangeable random partition of $[n]$ that is obtained when one uses the Poisson–Dirichlet $\text{PD}(\theta)$ prior to mix paintbox processes. Because of the equivalence between patterns and exchangeable random partitions, it

estimates the probability of a pattern ψ_1^n via the following formula:

$$q_\theta^{\text{CRP}}(\psi_1, \dots, \psi_n) = \frac{\theta^m}{\theta(\theta+1) \cdots (\theta+n-1)} \prod_{\mu=1}^n [(\mu-1)!]^{\phi_\mu}. \quad (1)$$

Recall that ϕ_μ is the number of symbols that appear μ times in ψ_1^n . In particular, the predictive distribution associated to the Ewens sampling formula or Chinese restaurant process is:

$$q_\theta^{\text{CRP}}(\psi|\psi_1, \dots, \psi_n) = \begin{cases} \frac{\mu}{n+\theta} & \psi \text{ appeared } \mu \text{ times} \\ & \text{in } \psi_1, \dots, \psi_n; \\ \frac{\theta}{n+\theta} & \psi \text{ is a new symbol.} \end{cases} \quad (2)$$

More generally, one can define the Pitman–Yor predictor (for $\alpha \in [0, 1]$ and $\theta > -\alpha$) as:

$$q_{\alpha, \theta}^{\text{PY}}(\psi|\psi_1, \dots, \psi_n) = \begin{cases} \frac{\mu-\alpha}{n+\theta} & \psi \text{ appeared } \mu \text{ times} \\ & \text{in } \psi_1, \dots, \psi_n; \\ \frac{\theta+m\alpha}{n+\theta} & \psi \text{ is a new symbol.} \end{cases} \quad (3)$$

where m is the number of distinct symbols in ψ_1^n . The probability assigned by the Pitman–Yor predictor to a pattern ψ_1^n is therefore:

$$q_{\alpha, \theta}^{\text{PY}}(\psi_1, \dots, \psi_n) = \frac{\theta(\theta+\alpha)(\theta+2\alpha) \cdots (\theta+(m-1)\alpha)}{\theta(\theta+1) \cdots (\theta+n-1)} \prod_{\mu=1}^n \left(\frac{\Gamma(\mu-\alpha)}{\Gamma(1-\alpha)} \right)^{\phi_\mu}. \quad (4)$$

Note that $\Gamma(\mu-\alpha)/\Gamma(1-\alpha) = (\mu-\alpha-1)(\mu-\alpha-2) \cdots (1-\alpha)$.

2.4 Strong Universality Measures: Worst-Case and Average

How should we measure the quality of a pattern probability predictor q ? We investigate two criteria here: the worst-case and the average-case redundancy. The redundancy of q on a given pattern ψ^n is:

$$R(q, \psi^n) \stackrel{\text{def}}{=} \sup_{p \in \mathcal{I}} \log \frac{p(\psi^n)}{q(\psi^n)},$$

The worst-case redundancy of q is defined to be:

$$\hat{R}(q) \stackrel{\text{def}}{=} \max_{\psi^n \in \Psi^n} \sup_{p \in \mathcal{I}} \log \frac{p(\psi^n)}{q(\psi^n)} = \sup_{p \in \mathcal{I}} \max_{\psi^n \in \Psi^n} \log \frac{p(\psi^n)}{q(\psi^n)}$$

Recall that $p(\psi^n)$ just denotes the probability that a length- n sequence generated by p has pattern ψ^n ; it is unnecessary to specify the support here. Since the set of length- n patterns is finite, there is no need for a supremum in the outer maximization above. The worst-case redundancy is often referred to as the per-sequence redundancy, as well.

The average-case redundancy replaces the max over patterns with an expectation over p :

$$\begin{aligned} \bar{R}(q) &\stackrel{\text{def}}{=} \sup_{p \in \mathcal{I}} \mathbb{E}_p \left[\log \frac{p(\psi^n)}{q(\psi^n)} \right] \\ &= \sup_{p \in \mathcal{I}} D(p \parallel q), \end{aligned}$$

where $D(\cdot\|\cdot)$ is the Kullback–Leibler divergence or relative entropy. That is, the average-case redundancy is nothing but the worst-case Kullback–Leibler divergence between the distribution p and the predictor q .

A pattern probability estimator is considered “good” if the worst-case or average-case redundancies are sublinear in n or $\hat{R}(q)/n \rightarrow 0$ and $\bar{R}(q)/n \rightarrow 0$ as $n \rightarrow \infty$. Succinctly put, redundancy that is sublinear in n implies that the underlying probability of a sequence can be estimated accurately almost surely. Redundancy is one way to measure the “frequentist” properties of the Bayesian approaches we consider in this paper and refers to the compressibility of the distribution from an information theoretic perspective.

As mentioned in the Introduction, redundancy differs from notions, such as the concentration of the posterior distribution about the true distribution. However, the parameters of the CRPs or the PY processes that compress well often correspond to the maximum likelihood (ML) estimates from the sample.

2.5 Weak Universality

In the previous section, we considered guarantees that hold over the entire model class; both the worst case and average case involve taking a supremum over the entire model class. Therefore, the strong guarantees—average or worst-case—hold uniformly over the model class. However, as we will see, exchangeable estimators, in particular the Chinese restaurant process and Pitman–Yor estimators, are tuned towards specific kinds of sources, rather than all i.i.d. models by the appropriate choice of parameters. This behavior is better captured by looking at the model-dependent convergence of the exchangeable estimators, which is known as weak universality. Specifically, let \mathcal{P}^∞ be a collection of i.i.d. measures over infinite sequences, and let \mathcal{P}_Ψ^∞ be the collection of measures induced on patterns by \mathcal{P}^∞ . We say an estimator q is weakly universal for a class \mathcal{P}_Ψ^∞ if for all $p \in \mathcal{P}^\infty$:

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \mathbb{E}_p \left[\log \frac{p(\psi^n)}{q(\psi^n)} \right] = 0.$$

3 Strong Redundancy

We now describe our main results on the redundancy of estimators derived from the prior distributions on \mathcal{I} .

3.1 Chinese Restaurant Process Predictors

Previously [41], it was shown by some of the authors that the worst-case and average-case redundancies for the CRP estimator are both $\Omega(n \log n)$, which means it is not strongly universal. However, this negative result follows because the CRP estimator is tuned not to the entire i.i.d. class of distributions, but to a specific subclass of i.i.d. sources depending on the choice of parameter. To investigate this further, we analyze the redundancy of the CRP estimator when we have a bound on the number m of distinct elements in the pattern ψ_1^n .

Chinese restaurant processes $q_\theta^{\text{CRP}}()$ with parameter θ are known to generate exchangeable random partitions where the number of distinct parts M satisfy $M/\log n \rightarrow \theta$ almost surely as the sample size n increases; see *e.g.*, [42]. Equivalently, the CRP generates patterns with M distinct symbols, where $M/\log n \rightarrow \theta$. However, the following theorem reverses the above setting. Here, we are given an i.i.d. sample of data of length n with m symbols (how the data were generated

is not important), but we pick the parameter of a CRP estimator that describes the pattern of the data well. While it is satisfying that the chosen parameter matches the ML estimate of the number of symbols, note that this need not necessarily be the only parameter choice that works.

Theorem 1. [Redundancy for CRP estimators] Consider the estimator $q_\theta^{\text{CRP}}(\psi_1^n)$ in (1) and (2). Then, for sufficiently large n and for patterns ψ_1^n whose number of distinct symbols m satisfies:

$$m \leq C \cdot \frac{n}{\log n} (\log \log n)^2,$$

the redundancy of the predictor $q_\theta^{\text{CRP}}(\psi_1^n)$ with $\lceil \theta \rceil = m / \log n$ satisfies:

$$\log \frac{p(\psi_1^n)}{q_\theta^{\text{CRP}}(\psi_1^n)} \leq 3C \cdot \frac{n(\log \log n)^3}{\log n} = o(n).$$

Proof The number of patterns with prevalence $\{\phi_\mu\}$ is:

$$\frac{n!}{\prod_{\mu=1}^n [\mu!]^{\phi_\mu} \phi_\mu!},$$

and therefore:

$$p(\psi_1^n) \leq \frac{\prod_{\mu=1}^n [\mu!]^{\phi_\mu} \phi_\mu!}{n!}, \quad (5)$$

since patterns with prevalence $\{\phi_\mu\}$ all have the same probability.

Using the upper bound in (5) on $p(\psi_1^n)$ and (1) yields:

$$\begin{aligned} \log \frac{p(\psi_1^n)}{q_\theta^{\text{CRP}}(\psi_1^n)} &\leq \log \prod_{\mu=1}^n \frac{[\mu!]^{\phi_\mu} \phi_\mu!}{[(\mu-1)!]^{\phi_\mu} \theta^m} + \log \frac{\theta(\theta+1) \cdots (\theta+n-1)}{n!} \\ &= \log \left(\prod_{\mu=1}^n \mu^{\phi_\mu} \right) + \log \left(\frac{1}{\theta^m} \prod_{\mu=1}^n \phi_\mu! \right) + \log \frac{\theta(\theta+1) \cdots (\theta+n-1)}{n!}. \end{aligned} \quad (6)$$

Let $\bar{\theta} = \lceil \theta \rceil$. The following bound follows from Stirling's approximation:

$$\frac{\theta(\theta+1) \cdots (\theta+n-1)}{n!} \leq \frac{(\bar{\theta}+n)!}{\theta! n!} \leq \frac{(\bar{\theta}+n)^{\bar{\theta}+n}}{\bar{\theta}^{\bar{\theta}} n^n} \leq \left(\frac{\bar{\theta}+n}{\bar{\theta}} \right)^{\bar{\theta}} \left(\frac{\bar{\theta}+n}{n} \right)^n. \quad (7)$$

The first term in (6) can be upper bounded by $\log(n/m)^m$ since the argument of the $\log(\cdot)$ is maximized over $\mu \in [1, n]$ when $\mu = n/m$. The second term is also maximized when all symbols appear the same number of times, corresponding to $\phi_\mu = m$ for one μ . Therefore:

$$\log \frac{p(\psi_1^n)}{q_\theta^{\text{CRP}}(\psi_1^n)} \leq \log \left(\frac{n}{m} \right)^m + \log \frac{m!}{\theta^m} + \log \left(\frac{\bar{\theta}+n}{\bar{\theta}} \right)^{\bar{\theta}} \left(1 + \frac{\bar{\theta}}{n} \right)^n.$$

Now, $\left(1 + \frac{\bar{\theta}}{n} \right)^n \leq e^{\bar{\theta}}$ for sufficiently large n , so:

$$\log \frac{p(\psi_1^n)}{q_\theta^{\text{CRP}}(\psi_1^n)} \leq \log \left(\frac{n}{m} \right)^m + \log \frac{m!}{\theta^m} + \log \left(\frac{(\bar{\theta}+n)e}{\bar{\theta}} \right)^{\bar{\theta}}. \quad (8)$$

Choose $\bar{\theta} = m/\log n$. This gives the bound:

$$\log \frac{p(\psi_1^n)}{q_{\bar{\theta}}^{\text{CRP}}(\psi_1^n)} \leq m \log \left(\frac{n}{m} \right) + \log \frac{m!}{m^m} \left(\frac{\bar{\theta}}{\theta} \right)^m + m \log \log n + \log \left(\frac{(\bar{\theta} + n)e}{\bar{\theta}} \right)^{\bar{\theta}}.$$

the second term is negative for sufficiently large m . Therefore:

$$\log \frac{p(\psi_1^n)}{q_{\bar{\theta}}^{\text{CRP}}(\psi_1^n)} \leq m \log \left(\frac{n}{m} \right) + m \log \log n + \frac{m}{\log n} \log \left(2 + \frac{n \log n}{m} \right). \quad (9)$$

Noting that the function above is monotonic in m for $n \geq 16$, we choose:

$$m = C \frac{n}{\log n} (\log \log n)^2,$$

and the bound becomes:

$$\begin{aligned} \log \frac{p(\psi_1^n)}{q_{\bar{\theta}}^{\text{CRP}}(\psi_1^n)} &\leq Cn \frac{(\log \log n)^2}{\log n} \log \left(\frac{\log n}{(\log \log n)^2} \right) + Cn \frac{(\log \log n)^3}{\log n} + Cn \left(\frac{\log \log n}{\log n} \right)^2 \log \left(2 + \left(\frac{\log n}{\log \log n} \right)^2 \right) \\ &\leq 3Cn \frac{(\log \log n)^3}{\log n} \\ &= o(n). \end{aligned} \quad \square$$

This theorem is slightly dissatisfying, since it requires us to have a bound on the number of distinct symbols in the pattern. In Section 4, we take mixtures of CRP estimators to arrive at estimators that are weakly universal.

3.2 Pitman–Yor Predictors

We now turn to the more general class of Pitman–Yor predictors. We can obtain a similar result as for the CRP estimator, but we can handle all patterns with $m = o(n)$ distinct symbols.

As before, the context for the following theorem is this: we are given an i.i.d. sample of data of length n with m symbols (there is no consideration, as before, as to how the data was generated), but we pick the parameters of a PY estimator that describes the pattern of the data well. The choice of the PY estimator is not necessarily the best, but one that will help us construct the weakly universal estimator in later sections of this paper.

We also note that the choice of the parameter θ below is analogous to our choice when $\alpha = 0$ (reducing to the CRP case). For patterns generated by a PY process, where $0 < \alpha \leq 1$, the number of distinct symbols grows as n^α . It is known that in this regime, the choice of θ is not distinguishable [13]. However, what is known is that the choice of θ remains $o(n^\alpha)$, something that is achieved in the selection of θ below. As the reader will note, as long as $0 < \alpha < 1$ is fixed, the theorem below will help us construct weakly universal estimators further on.

Theorem 2. [Worst-case redundancy] Consider the estimator $q_{\alpha, \theta}^{\text{PY}}(\psi_1^n)$. Then, for sufficiently large n and for patterns ψ_1^n , whose number of distinct symbols m satisfies $m = o(n)$, the worst-case redundancy of the predictor $q_{\alpha, \theta}^{\text{PY}}(\psi_1^n)$ with $\theta = m/\log n$ satisfies:

$$\log \frac{p(\psi_1^n)}{q_{\alpha, \theta}^{\text{PY}}(\psi_1^n)} = o(n). \quad (10)$$

Proof For a pattern ψ_1^n , from the definition of $q_{\alpha,\theta}^{\text{PY}}(\psi_1^n)$ in (4) and (5),

$$\log \frac{p(\psi_1^n)}{q_{\alpha,\theta}^{\text{PY}}(\psi_1^n)} \leq \log \left(\frac{\prod_{\mu=1}^n [\mu!]^{\phi_\mu} \phi_\mu!}{n!} \cdot \frac{(\theta+1) \cdots (\theta+n-1)}{(\theta+\alpha)(\theta+2\alpha) \cdots (\theta+m\alpha)} \prod_{\mu=1}^n \left(\frac{\Gamma(1-\alpha)}{\Gamma(\mu-\alpha-1)} \right)^{\phi_\mu} \right). \quad (11)$$

We can bound the components separately. First, as before we have:

$$\prod_{\mu=1}^n \phi_\mu! \leq m!$$

Since $\theta > -\alpha$, we have $\theta + \alpha > 0$ and:

$$\begin{aligned} (\theta + \alpha)(\theta + 2\alpha) \cdots (\theta + (m-1)\alpha) &\geq (\theta + \alpha)\alpha(2\alpha) \cdots ((m-2)\alpha) \\ &= (\theta + \alpha)(m-2)!\alpha^{m-2}. \end{aligned}$$

Again, letting $\bar{\theta} = \lceil \theta \rceil$, from the same arguments as in (7) and (8),

$$\log \frac{\theta(\theta+1) \cdots (\theta+n-1)}{n!} \leq \bar{\theta} \log \frac{(\bar{\theta}+n)e}{\bar{\theta}}.$$

Finally, note that $(1-\alpha)(2-\alpha) \cdots (\mu-1-\alpha) \geq (1-\alpha)(\mu-2)!$, so:

$$\begin{aligned} \frac{\prod_{\mu=1}^n [\mu!]^{\phi_\mu}}{\prod_{\mu=1}^n [(1-\alpha)(2-\alpha) \cdots (\mu-1-\alpha)]^{\phi_\mu}} &\leq \prod_{\mu=1}^n \left(\frac{\mu!}{(1-\alpha)(\mu-2)!} \right)^{\phi_\mu} \\ &\leq \frac{\prod_{\mu=1}^n \mu^{2\phi_\mu}}{(1-\alpha)^m} \\ &\leq \frac{(n/m)^{2m}}{(1-\alpha)^m}. \end{aligned}$$

Putting this together:

$$\begin{aligned} \log \frac{p(\psi_1^n)}{q_{\alpha,\theta}^{\text{PY}}(\psi_1^n)} &\leq \log \frac{m!}{(\theta + \alpha)(m-2)!\alpha^{m-2}} + \bar{\theta} \log \frac{(\bar{\theta}+n)e}{\bar{\theta}} + \log \frac{(n/m)^{2m}}{(1-\alpha)^m} \\ &\leq 2m \log \frac{n}{m} + (m-2) \log \frac{1}{(1-\alpha)\alpha} + \bar{\theta} \log \frac{(\bar{\theta}+n)e}{\bar{\theta}} + \log \frac{m^2}{(\theta + \alpha)} + \log \frac{1}{(1-\alpha)^2}. \end{aligned} \quad (12)$$

If $m = o(n)$, then the right side above is less than $o(n)$, as desired. \square

It is well known that the Pitman–Yor process can produce patterns whose relative frequency is zero, e.g. the pattern $1^k 23 \cdots (n-k)$. Therefore, it is not surprising that the worst-case redundancy and average case redundancies can be bad. However, as the next theorem shows, the actual redundancy of the Pitman–Yor estimator is $\Theta(n)$, which is significantly better than the lower bound of $\Omega(n \log n)$ proven in Santhanam and Madiman [41] for Chinese restaurant processes.

Theorem 3. [Redundancies] Consider the estimator $q_{\alpha,\theta}^{\text{PY}}(\psi_1^n)$. Then, for sufficiently large n , the worst-case redundancy and average case redundancy satisfy:

$$\hat{R}(q_{\alpha,\theta}^{\text{PY}}(\cdot)) = \Theta(n) \quad \text{and} \quad \bar{R}(q_{\alpha,\theta}^{\text{PY}}(\cdot)) = \Theta(n). \quad (13)$$

That is, $q_{\alpha,\theta}^{\text{PY}}(\cdot)$ is neither strongly nor weakly universal.

Proof For the upper bound, we start with (12) and note that in the worst case, $m = O(n)$, so $\hat{R}(q_{\alpha,\theta}^{\text{PY}}(\cdot)) = O(n)$ and *a fortiori* $\bar{R}(q_{\alpha,\theta}^{\text{PY}}(\cdot)) = O(n)$.

For the lower bound, consider the patterns $11 \cdots 1$ and $12 \cdots n$. For the Pitman–Yor estimator,

$$\begin{aligned} q_{\alpha,\theta}^{\text{PY}}(11 \cdots 1) q_{\alpha,\theta}^{\text{PY}}(12 \cdots n) &= \frac{\theta(1-\alpha) \cdots (n-1-\alpha)}{\theta(\theta+1) \cdots (\theta+n-1)} \frac{\theta(\theta+\alpha) \cdots (\theta+(n-1)\alpha)}{\theta(\theta+1) \cdots (\theta+n-1)} \\ &= \frac{(1-\alpha)(\theta+\alpha)}{(\theta+1)^2} \cdot \frac{(2-\alpha)(\theta+2\alpha)}{(\theta+2)^2} \cdots \frac{(n-1-\alpha)(\theta+(n-1)\alpha)}{(\theta+n-1)^2}. \end{aligned}$$

For $j \geq 1$, $0 < \alpha < 1$ and $\alpha + \theta > 0$, we show in the claim proven below that:

$$\frac{(j-\alpha)(\theta+j\alpha)}{(\theta+j)^2} \leq \max\left\{\frac{1}{2}, \alpha\right\}.$$

Therefore, each term is less than one. Then, for $\alpha < 1$, there exists a constant $0 < c < 1$, such that:

$$q_{\alpha,\theta}^{\text{PY}}(11 \cdots 1) q_{\alpha,\theta}^{\text{PY}}(12 \cdots n) \leq c^n.$$

Thus:

$$\log \frac{1}{q_{\alpha,\theta}^{\text{PY}}(11 \cdots 1)} + \log \frac{1}{q_{\alpha,\theta}^{\text{PY}}(12 \cdots n)} \geq n \log \frac{1}{c}.$$

Let the distribution p_1 be a singleton, so $p_1(1 \cdots 1) = 1$. For any small $\delta > 0$, we can find a distribution p_n , such that $p_n(12 \cdots n) = 1 - \delta$ by choosing p_n to be uniform on a sufficiently large set. Thus:

$$\begin{aligned} \hat{R}(q_{\alpha,\theta}^{\text{PY}}(\cdot)) &\geq \max \left\{ \log \frac{1-\delta}{q_{\alpha,\theta}^{\text{PY}}(11 \cdots 1)}, \log \frac{1-\delta}{q_{\alpha,\theta}^{\text{PY}}(12 \cdots n)} \right\} \\ &\geq \frac{1}{2} \left(\log \frac{1}{q_{\alpha,\theta}^{\text{PY}}(11 \cdots 1)} + \log \frac{1}{q_{\alpha,\theta}^{\text{PY}}(12 \cdots n)} \right) + \log(1-\delta) \\ &\geq \frac{n}{2} \log \frac{1}{c} + \log(1-\delta). \end{aligned}$$

This shows that $\hat{R}(q_{\alpha,\theta}^{\text{PY}}(\cdot)) = \Omega(n)$. Furthermore,

$$\begin{aligned} \bar{R}(q_{\alpha,\theta}^{\text{PY}}(\cdot)) &\geq \max \left\{ (1-\delta) \frac{1-\delta}{q_{\alpha,\theta}^{\text{PY}}(11 \cdots 1)}, (1-\delta) \frac{1-\delta}{q_{\alpha,\theta}^{\text{PY}}(12 \cdots n)} \right\} \\ &\geq (1-\delta) \left(\frac{n}{2} \log \frac{1}{c} + \log(1-\delta) \right), \end{aligned}$$

so $\bar{R}(q_{\alpha,\theta}^{\text{PY}}(\cdot)) = \Omega(n)$.

All that remains is to prove the following claim:

Claim 1. For $j \geq 1$, $0 < \alpha < 1$ and $\alpha + \theta > 0$ we show that:

$$\frac{(j - \alpha)(\theta + j\alpha)}{(\theta + j)^2} \leq \max\{\frac{1}{2}, \alpha\}.$$

Proof First, assume that $0 < \alpha < \frac{1}{2}$. Then, the inequality is:

$$\frac{(j - \alpha)(\theta + j\alpha)}{(\theta + j)^2} \leq \frac{1}{2}.$$

Equivalently, we need to show:

$$(1 - 2\alpha)j^2 + 2\alpha^2j + \theta^2 + 2\alpha\theta \geq 0.$$

Since $1 - 2\alpha > 0$, the quadratic is always nondecreasing on $j \geq 1$. Therefore, the positive integer $j = 1$ minimizes the quadratic above, and the value of the quadratic at $j = 1$ is:

$$1 - 2\alpha + 2\alpha^2 + \theta^2 + 2\alpha\theta = (\alpha - 1)^2 + (\theta + \alpha)^2 \geq 0.$$

Next, assume that $\frac{1}{2} \leq \alpha < 1$. Then, the inequality is:

$$\frac{(j - \alpha)(\theta + j\alpha)}{(\theta + j)^2} \leq \alpha.$$

Equivalently, we need to show:

$$((2\alpha - 1)\theta + \alpha^2)j + \alpha\theta(\theta + 1) \geq 0. \tag{14}$$

Since $2\alpha - 1 \geq 0$ and $\theta > -\alpha$,

$$(2\alpha - 1)\theta + \alpha^2 \geq -(2\alpha - 1)\alpha + \alpha^2 = \alpha(1 - \alpha) > 0.$$

Therefore, the minimum of the left equation in (14) is achieved at $j = 1$. Note that $\alpha\theta^2 > -\alpha^3 - 2\alpha^2\theta$ by using $(\alpha + \theta)^2 > 0$. Therefore, the value of the minimum is:

$$\begin{aligned} (2\alpha - 1)\theta + \alpha^2 + \alpha\theta(\theta + 1) &= \alpha\theta^2 + (3\alpha - 1)\theta + \alpha^2 \geq -\alpha^3 + (-2\alpha^2 + 3\alpha - 1)\theta + \alpha^2 \\ &\geq -\alpha^3 - (-2\alpha^2 + 3\alpha - 1)\alpha + \alpha^2 \\ &= \alpha^3 - 2\alpha^2 + \alpha = \alpha(\alpha - 1)^2 \geq 0. \end{aligned}$$

Note that $-2\alpha^2 + 3\alpha - 1 \geq 0$ for $\frac{1}{2} \leq \alpha < 1$, and the claim follows. \square

The theorem follows. \square

4 Weak Universality

In this section, we show how to modify the CRP or PY estimators to obtain weakly universal estimators. The CRP and PY cases are identical; therefore, we only work out the CRP case.

For all $i \geq 1$ and $j \geq 1$, let:

$$c_{i,j} = \frac{1}{i(i+1)j(j+1)}$$

so that $\sum_{i,j} c_{i,j} = 1$. Let $\tilde{q}_{i,j}^{\text{CRP}}$ be the CRP measure over patterns with $\theta = i/\log j$. Consider the following measure over patterns of infinite sequences that assigns, for all n and all patterns ψ^n of length n , the probability:

$$q^*(\psi^n) = \sum_{i,j} c_{i,j} \tilde{q}_{i,j}^{\text{CRP}} \psi^n. \quad (15)$$

We will show that q^* is a weakly universal measure over patterns of i.i.d. sequences.

To do so, we will need the following two lemmas. Lemma 4 is a useful “folk” inequality that we believe is attributed to Minkowski. Lemma 5 relates the expected number of distinct symbols in length n sequences of an i.i.d. process to its entropy and is of independent interest. The result not only strengthens a similar result in [43], but also provides a different and more compact proof.

Lemma 4. For $n \geq 1$, let $x_1 \geq x_2 \geq \dots x_n \geq 0$ and $y_1 \geq y_2 \geq \dots y_n \geq 0$ be two sorted sequences. Then:

$$\frac{1}{n} \sum_{l=1}^n x_l y_l \geq \left(\frac{1}{n} \sum_{i=1}^n x_i \right) \left(\frac{1}{n} \sum_{j=1}^n y_j \right) \geq \frac{1}{n} \sum_{l=1}^n x_l y_{n+1-l}.$$

Proof The left inequality of the lemma follows by noting that:

$$\left(\sum_{i=1}^n x_i \right) \left(\sum_{j=1}^n y_j \right) = \sum_{k=0}^{n-1} \sum_{l=1}^n x_l y_{l+k},$$

and that the sum $\sum_l x_l y_{l+k}$ is maximized at $k = 0$, since both sequences are sorted in the same direction. The right inequality of the lemma can be proven similarly, but will not be used in the paper. \square

Lemma 5. For all discrete i.i.d. processes P with entropy rate (or marginal entropy) H , let M_n be the random variable counting the number of distinct symbols in a sample of length n drawn from P . The following bound holds:

$$\mathbb{E}[M_n] \leq \frac{nH}{\log n} + 1. \quad (16)$$

Proof Let $P(i) = p_i$. We begin by noting that:

$$H = \sum_i p_i \log \frac{1}{p_i} = \sum_i p_i \sum_{j=1}^{\infty} \frac{(1-p_i)^j}{j},$$

where the second equality follows by the Taylor series expansion:

$$-\log p_i = -\log(1 - (1-p_i)) = \sum_{j=1}^{\infty} \frac{(1-p_i)^j}{j}.$$

The right summation in the equation above is bounded below as follows:

$$\begin{aligned} \sum_{j=1}^{\infty} \frac{(1-p_i)^j}{j} &\geq \sum_{j=1}^n \frac{(1-p_i)^j}{j} \\ &\stackrel{(a)}{\geq} \frac{1}{n} \sum_{l=1}^n \frac{1}{l} \sum_{m=1}^n (1-p_i)^m \\ &\geq \frac{\log n}{n} \frac{(1-p_i)}{p_i} (1 - (1-p_i)^n) \end{aligned}$$

where (a) follows from Minkowski's inequality in Lemma 4 and the last inequality, because $\sum_{l=1}^n \frac{1}{l} \geq \log n$. Thus,

$$H = \sum_i p_i \sum_{j=1}^{\infty} \frac{(1-p_i)^j}{j} \geq \frac{\log n}{n} \sum_i (1-p_i)(1-(1-p_i)^n) \geq \frac{\log n}{n}(\mathbb{E}M_n - 1)$$

where for the second inequality, we use $\sum_i p_i(1-(1-p_i)^n) \leq \sum_i p_i \leq 1$. \square

Theorem 6. [Weak universality for CRP mixtures] For all discrete i.i.d. processes $p \in \mathcal{I}$ with a finite entropy rate,

$$D(p \parallel q^*) = o(n). \quad (17)$$

That is, q^* is weakly universal.

Proof We write the divergence between p and q^* in (15) as the expected log ratio and condition on the value of M_n :

$$\begin{aligned} & \mathbb{E}_p \left[\log \frac{p(\psi_1^n)}{\sum_m c_{m,n} \tilde{q}_{m,n}^{\text{CRP}}(\psi_1^n)} \right] \\ = & \mathbb{P} \left(M_n > \frac{n(\log \log n)^2}{\log n} \right) \cdot \mathbb{E}_p \left[\log \frac{p(\psi_1^n)}{\sum_m c_{m,n} \tilde{q}_{m,n}^{\text{CRP}}(\psi_1^n)} \middle| M_n > \frac{n(\log \log n)^2}{\log n} \right] \\ & + \mathbb{P} \left(M_n < \frac{n(\log \log n)^2}{\log n} \right) \cdot \mathbb{E}_p \left[\log \frac{p(\psi_1^n)}{\sum_m c_{m,n} \tilde{q}_{m,n}^{\text{CRP}}(\psi_1^n)} \middle| M_n < \frac{n(\log \log n)^2}{\log n} \right]. \end{aligned} \quad (18)$$

Consider the estimator $\tilde{q}_{i,j}^{\text{CRP}} \psi_1^n$ in q^* corresponding to $i = M_n$ and $j = \log n$. This is the estimator $q_{\theta}^{\text{CRP}}(\psi_1^n)$ with $\theta = M_n/\log n$. From the proof of Theorem 1, we have:

$$\begin{aligned} \log \frac{p(\psi_1^n)}{\sum_m c_{m,n} \tilde{q}_{m,n}^{\text{CRP}}(\psi_1^n)} & \leq \log \frac{p(\psi_1^n)}{c_{i,j} \tilde{q}_{i,j}^{\text{CRP}} \psi_1^n} \\ & \leq \log \frac{1}{c_{ij}} + \log \frac{p(\psi_1^n)}{\tilde{q}_{i,j}^{\text{CRP}} \psi_1^n} \\ & \leq \log (M_n(M_n + 1)(\log n)(\log n + 1)) + \log \frac{p(\psi_1^n)}{q_{\theta}^{\text{CRP}}(\psi_1^n)} \end{aligned} \quad (19)$$

We will bound the two terms in (19) in the regimes for M_n .

The result of Theorem 1 says that if $M_n < \frac{n(\log \log n)^2}{\log n}$, then:

$$\log \frac{1}{c_{ij}} + \log \frac{p(\psi_1^n)}{q_{\theta}^{\text{CRP}}(\psi_1^n)} \leq \log (M_n(M_n + 1)(\log n)(\log n + 1)) + o(n) \quad (20)$$

Thus, this term is $o(n)$.

If $M_n > \frac{n(\log \log n)^2}{\log n}$, then we first apply Markov's inequality using the previous lemma:

$$\begin{aligned} \mathbb{P} \left(M_n > \frac{n(\log \log n)^2}{\log n} \right) & \leq \frac{\log n}{n(\log \log n)^2} \left(\frac{nH}{\log n} + 1 \right) \\ & \leq \frac{H}{(\log \log n)^2} + \frac{\log n}{n(\log \log n)^2}. \end{aligned} \quad (21)$$

Therefore, for all finite entropy processes, this probability goes to zero as $n \rightarrow \infty$. Looking at the term in q^* corresponding $q_\theta^{\text{CRP}}(\psi_1^n)$ with $\theta = M_n/\log n$ and using the fact that $M_n \leq n$, we see that the first term in (19) is upper bounded as $O(\log n)$. For the second term, we appeal to (9) in the proof of Theorem 1:

$$\log \frac{p(\psi_1^n)}{q_\theta^{\text{CRP}}(\psi_1^n)} \leq M_n \log \frac{M_n}{n} + M_n \log \log n + \frac{M_n}{\log n} \log \left(2 + \frac{n \log n}{M_n} \right) \quad (22)$$

$$\leq n + n \log \log n + n \frac{\log(2 + n \log n)}{\log n} \quad (23)$$

$$\leq 3n \log \log n. \quad (24)$$

Plugging these terms into (18):

$$D(p \parallel q^*) \leq \left(\frac{H}{(\log \log n)^2} + \frac{\log n}{n(\log \log n)^2} \right) \cdot O(n \log \log n) + 1 \cdot o(n) = o(n). \quad \square$$

The preceding theorem shows that the mixture of CRP estimators q^* is weakly universal. However, note that q^* is not itself a CRP estimator. An identical construction is possible for the PY estimators, as well. The convergence of the weakly universal q^* depends on the number of entropy of the source, as well as the number of distinct symbols in a sample of size n .

While it would be tempting to predict the performance of the estimator q^* for larger sample sizes $N \geq n$, such a task requires a more careful analysis. In general, it may be impossible to non-trivially bound the number of distinct symbols M_N with smaller sample size n , as the following example shows.

Example 1. Let $n = \sqrt{N}$. Consider a set \mathcal{I} containing the following two distributions: (i) p over \mathbb{N} , which assigns probability $1 - 1/n^{3/2} = 1 - 1/N^{3/4}$ to the atom 1 and splitting the probability $1/N^{3/4}$ equally among the elements of the set $\{2, \dots, N^2 + 1\}$; and (ii) p' , which simply assigns probability one to one. A sample of size n from either p or p' is 1^n with probability at least $1 - 1/N^{1/4}$, no matter what the underlying source is; therefore, we cannot distinguish between these sources with probability $1 - 1/N^{1/4}$ from a sample of size n .

However, a sample of size N from p has $\mathcal{O}(N^{1/4})$ distinct symbols on average, while that of p' will have only one element. It follows that if all we know is that the unknown distribution comes from \mathcal{I} , with high probability under the unknown source, we cannot predict whether the number of symbols in a sample of size N will remain one or not from a sample of length n . Furthermore, by changing the ratio of n and N (and therefore, the probability of the symbol 1 under p), we can make the expected number of symbols in a N -length sample under p as large as we want. \square

However, it is possible to impose restrictions on the class of distributions that allow us to ensure that we can predict the number of symbols in longer samples. In future work, we will borrow from the data-derived consistency formulations of [44] to characterize when we will be able to predict the number of symbols in longer samples.

5 Conclusions and Future Work

In this note, we investigated the worst-case and average-case redundancies of pattern probability estimators derived from priors on \mathcal{I} that are popular in Bayesian statistics. Both the CRP and Pitman–Yor estimators give a vanishing redundancy per symbol for patterns whose number of

distinct symbols m is sufficiently small. The Pitman–Yor estimator requires only that $m = o(n)$, which is an improvement on the CRP. However, when m can be arbitrarily large (or the alphabet size is arbitrarily large), the worst-case and average-case redundancies do not scale like $o(n)$. Here, again, the Pitman–Yor estimator is superior, in that the redundancies scale like $\Theta(n)$ as opposed to the $\Omega(n \log n)$ for the CRP estimator. While these results show that these estimators are not strongly universal, we constructed a mixture of CRP process (which is not itself a CRP estimator) that is weakly universal. One of the estimators derived in [16] is exchangeable and has near-optimal worst case redundancy of $O(\sqrt{n})$. Kingman’s results imply this estimator corresponds to a prior on \mathcal{I} ; however, this prior is yet unknown. Finding this prior may potentially reveal new interesting classes of priors other than the Poisson–Dirichlet priors.

Acknowledgments

The authors thank the American Institute of Mathematics and NSF for sponsoring a workshop on probability estimation, as well as A. Orlitsky and K. Viswanathan, who co-organized the workshop with the first author. They additionally thank P. Diaconis, M. Dudik, F. Chung, R. Graham, S. Holmes, O. Milenkovic, O. Shayevitz, A. Wagner, J. Zhang and M. Madiman for helpful discussions.

N.P. Santhanam was supported by a startup grant from the University of Hawaii and NSF Grants CCF-1018984 and CCF-1065632. A.D. Sarwate was supported in part by the California Institute for Telecommunications and Information Technology (CALIT2) at the University of California, San Diego. J.O. Woo was supported by NSF Grant CCF-1065494 and CCF-1346564.

References

- [1] P. Laplace, *Philosophical essay on probabilities. Translated from the fifth French edition of 1825*. No. 13 in Sources in the History of Mathematics and Physical Sciences, Springer Verlag, New York, 1995.
- [2] A. De Morgan, *An Essay on Probabilities, and on their Application to Life Contingencies and Insurance Offices*,. London, UK: Longman, Orme, Brown, Green & Longmans, 1838.
- [3] A. De Morgan, “Theory of probabilities,” in *Encyclopædia Metropolitana* (E. Smedley, H. J. Rose, and H. J. Rose, eds.), vol. II (Pure Sciences), pp. 393–490, B. Fellowes et. al., London, 1845.
- [4] I. Good, “The population frequencies of species and the estimation of population parameters,” *Biometrika*, vol. 40, pp. 237–264, December 1953.
- [5] D. Blackwell and J. B. MacQueen, “Ferguson distributions via Pólya urn schemes,” *The Annals of Statistics*, vol. 1, no. 2, pp. 353–355, 1973.
- [6] J. F. C. Kingman, “Random discrete distributions,” *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 37, no. 1, pp. 1–22, 1975.
- [7] J. Kingman, “Random partitions in population genetics,” *Proceedings of the Royal Society*, vol. 361, pp. 1–20, May 1978.
- [8] J. Kingman, “The representation of partition structures,” *Journal of the London Mathematical Society*, vol. s2-18, pp. 374–380, October 1978.

- [9] P. Diaconis and D. Freedman, “De Finetti’s generalizations of exchangeability,” in *Studies in Inductive Logic and Probability* (R. C. Jeffrey, ed.), vol. 2, pp. 233–250, Berkeley and Los Angeles, CA, USA: University of California Press, 1980.
- [10] D. J. Aldous, “Exchangeability and related topics,” in *École d’été de probabilités de Saint-Flour, XIII—1983*, no. 1117 in Lecture Notes in Mathematics, pp. 1–198, Berlin: Springer, 1985.
- [11] S. Zabell, “Predicting the unpredictable,” *Synthese*, vol. 90, no. 2, pp. 205–232, 1992.
- [12] J. Pitman, “Exchangeable and partially exchangeable random partitions,” *Probability Theory and Related Fields*, vol. 102, no. 2, pp. 145–158, 1995.
- [13] J. Pitman and M. Yor, “The two-parameter Poisson-Dirichlet distribution derived from a stable subordinator,” *Annals of Probability*, vol. 25, no. 2, pp. 855–900, 1997.
- [14] B. Clarke and A. Barron, “Information theoretic asymptotics of Bayes methods,” *IEEE Transactions on Information Theory*, vol. 36, pp. 453–471, May 1990.
- [15] B. Clarke and A. Barron, “Jeffreys’ prior is asymptotically least favorable under entropy risk,” *Journal of Statistical Planning and Inference*, vol. 41, no. 1, pp. 37–60, 1994.
- [16] A. Orlitsky, N. Santhanam, and J. Zhang, “Universal compression of memoryless sources over unknown alphabets,” *IEEE Transactions on Information Theory*, vol. 50, pp. 1469–1481, July 2004.
- [17] A. Orlitsky, N. Santhanam, and J. Zhang, “Always Good Turing: Asymptotically optimal probability estimation,” *Science*, vol. 302, pp. 427–431, October 17 2003. See also *Proceedings of the 44th Annual Symposium on Foundations of Computer Science*, October 2003.
- [18] B. Ryabko, “Compression based methods for nonparametric on-line prediction, regression, classification and density estimation of time series,” in *Festschrift in Honor of Jorma Rissanen on the occasion of his 75th birthday* (P. Grünwald, P. Myllymäki, I. Tabus, M. Weinberger, and B. Yu, eds.), pp. 271–288, Tampere International Center for Signal Processing, 2008.
- [19] A. B. Wagner, P. Viswanath, and S. R. Kulkarni, “Probability estimation in the rare-events regime,” *IEEE Transactions on Information Theory*, vol. 57, pp. 3207–3229, June 2011.
- [20] A. Nadas, “Good, Jelinek, Mercer, and Robins on Turing’s estimate of probabilities,” *American Journal of Mathematical and Management Sciences*, vol. 11, pp. 229–308, 1991.
- [21] W. Gale and K. Church, “What is wrong with adding one?,” in *Corpus based research into language* (N. Oostdijk and P. de Haan, eds.), pp. 189–198, Amsterdam, The Netherlands: Rodopi, 1994.
- [22] D. McAllester and R. Schapire, “On the convergence rate of Good Turing estimators,” in *Proceedings of the Thirteenth Annual Conference on Computational Learning Theory (COLT 2000)*, June 28 - July 1, 2000, Palo Alto, California (N. Cesa-Bianchi and S. A. Goldman, eds.), (San Francisco, CA), pp. 1–6, Morgan Kaufmann, June–July 2000.
- [23] E. Drukh and Y. Mansour, “Concentration bounds on unigrams language models,” *Journal of Machine Learning Research*, vol. 6, pp. 1231–1264, August 2005.
- [24] A. Orlitsky, N. Santhanam, K. Viswanathan, and J. Zhang, “On modeling profiles instead of values,” in *Proceedings of the Twentieth Conference Conference on Uncertainty in Artificial Intelligence (UAI-04)* (C. Meek and J. Halpern, eds.), (Arlington, VA, USA), pp. 426–435, AUAI Press, 2004.

- [25] R. Gallager, “Source coding with side information and universal coding,” Tech. Rep. LIDS-P-937, M.I.T., Cambridge, MA, USA, Sep 1976.
- [26] L. Davisson and A. Leon-Garcia, “A source matching approach to finding minimax codes,” *IEEE Transactions on Information Theory*, vol. 26, pp. 166–174, March 1980.
- [27] B. Ya. Ryabko, “Coding of a source with unknown but ordered probabilities,” *Problems of Information Transmission*, vol. 15, pp. 134–138, Oct 1979.
- [28] J. Kingman, *The Mathematics of Genetic Diversity*. No. 34 in CBMS-NSF Regional Conference Series in Applied Mathematics, Philadelphia, PA, USA: Society for Industrial and Applied Mathematics, 1980.
- [29] S. Zabell, “The continuum of inductive methods revisited,” in *The Cosmos of Science: Essays of Exploration* (J. Earman and J. D. Norton, eds.), ch. 12, Pittsburgh, PA, USA: The University of Pittsburgh Press, 1997.
- [30] S. Zabell, *Symmetry and Its Discontents: Essays on the History of Inductive Probability*. Cambridge Studies in Probability, Induction, and Decision Theory, Cambridge: Cambridge University Press, 2005.
- [31] J. Pitman, “Random discrete distributions invariant under size-biased permutation,” *Advances in Applied Probability*, vol. 28, pp. 525–539, 1996.
- [32] J. Acharya, H. Das, A. Jafarpour, A. Orlitsky, and A. Suresh, “Tight bounds for universal compression of large alphabets,” in *Proceedings of the IEEE International Symposium on Information Theory (ISIT)*, 2013.
- [33] J. Acharya, H. Das, and A. Orlitsky, “Tight bounds on profile redundancy and distinguishability,” in *Advances in Neural Information Processing Systems 25* (F. Pereira, C. Burges, L. Bottou, and K. Weinberger, eds.), pp. 3257–3265, Curran Associates, Inc., 2012.
- [34] J. Acharya, A. Jafarpour, A. Orlitsky, and A. Suresh, “Optimal probability estimation with applications to prediction and classification,” in *Conference on Learning Theory, vol. 30 of JMLR Workshop & Conference Proceedings*, 2013.
- [35] J. Pitman, *Combinatorial stochastic processes*, vol. 1875 of *Lecture Notes in Mathematics*. Berlin: Springer-Verlag, 2006.
- [36] T. Ferguson, “A Bayesian analysis of some nonparametric problems,” *The Annals of Statistics*, vol. 1, no. 2, pp. 209–230, 1973.
- [37] R. Ramamoorthi and K. Srikanth, “Dirichlet processes,” in *Encyclopedia of Statistical Sciences*, New York, NY, USA: John Wiley and Sons, 2007.
- [38] W. J. Ewens, “The sampling theory of selectively neutral alleles,” *Theoretical Population Biology*, vol. 3, pp. 87–112, March 1972.
- [39] S. Karlin and J. McGregor, “Addendum to a paper of W. Ewens,” *Theoretical Population Biology*, vol. 3, pp. 113–116, March 1972.
- [40] G. A. Watterson, “The sampling theory of selectively neutral alleles,” *Advances in Applied Probability*, vol. 6, pp. 463–488, September 1974.
- [41] N. Santhanam and M. Madiman, “Patterns and exchangeability,” in *2010 IEEE International Symposium on Information Theory (ISIT)*, (Austin, Texas, USA), pp. 1483–1487, June 2010.
- [42] M. A. Carlton, *Applications of the two-parameter Poisson-Dirichlet distribution*. PhD thesis, UNIVERSITY OF CALIFORNIA Los Angeles, 1999.

- [43] A. Orlitsky, N. P. Santhanam, K. Viswanathan, and J. Zhang, “Limit results on pattern entropy,” *Information Theory, IEEE Transactions on*, vol. 52, no. 7, pp. 2954–2964, 2006.
- [44] N. Santhanam, V. Anantharam, A. Kavcic, and W. Szpankowski, “Data driven weak universal redundancy,” in *Proceedings of IEEE Symposium on Information Theory*, Jul 2014.